

SEQUENCE MATCHING, SIMPLE SEARCHING

PGA Course in Bioinformatics
Tools for Comparative Analysis
May 12, 2003

Outline

- Sequence alignment algorithms
 - Rigorous Optimality: Needleman-Wunsch and Smith-Waterman
 - Rapid, heuristic algorithms
 - BLAST
 - FASTA
 - and their relatives
- Databases and Search Tools

What do you want to do?

- See if it is already in a database?
- Find similar sequences?
- Identify homologous sequences?
- Consider functional assignments and other annotation?
- Find primers in sequence?
- Find a short peptide?

UTILITY

- Relationships in evolution
- Identification of genes
- Assignment of possible functionality for genes or residues
- Possible structural understanding
- Aid in sequence assembly

What are you Comparing

• **Homologue**

Sequences that share a common ancestor; may have similar function

• **Paralogue**

Similar sequence within species, may have similar function

• **Orthologue**

Same sequence separated by a speciation event, probably same function

ANALOG

Non-homologue proteins that have similar folding architecture, or similar functional sites, which are believed to have arisen through convergent evolution

Searching for Similarity

• BLAST

- Search at NCBI and other servers (or locally)
- Non-redundant set of databases, one DB at a time
- Fast
- Shows several similar regions

Searching for Similarity

• FASTA

- Search against user-defined search sets, DB or subsections
- Only the single most similar region is shown

The Word –Size Parameter

A word is any short sequence less than or equal to six letter

- Protein 1-2
- Nucleotide 1-6

High word Size

- Faster
- Less Sensitive
- More Selective

Evolution and Alignment

Evolutionary concepts enable the determination of similarity and homology

- Similarity is an observable quantity, such as %identity
- Homology is a conclusion drawn from the data that two genes share a common evolutionary history.

Evolution and Alignments (2)

- ☛ Genes are either homologous or not homologous.
- ☛ There is no degree of homology
- ☛ You can't tell what the ancestral sequence is simply because you have two or more homologues.

So, what IS an Alignment?

Evolution and Alignments (3)

- ☛ Alignments reflect the PROBABLE evolutionary history of two sequences
- ☛ Residues that align and are not identical represent substitutions
- ☛ Sequences without correspondence in aligned sequences are interpreted as indels and in an alignment are gaps.

Evolution and Alignment

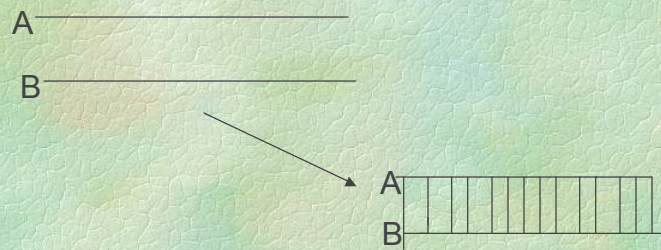
- Certain regions are more conserved than others, based on structure/function
- Certain regions may be conserved simply by history, not function
- This is true especially for closely related species.

Structure and Alignment

- If two proteins have more than 20-30% ID aligned, then the 3-D structures tend to be similar
- Overall folds are the same, details differ
- Form often follows function (Beware the BUT).
- So, sequence alignment is sometimes a 3-D alignment.

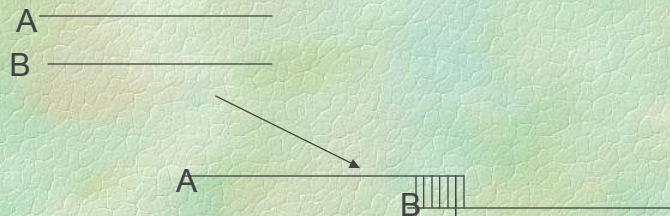
Global Alignment

Optimal alignment over the entire length



Local Alignment

Finds the highest scoring alignment regardless of position and length



Needleman Wunsch Algorithm

- Global alignment:: every residue of the two sequences has to participate
- Guaranteed to calculate an Optimal similarity score
- Begin at the beginning of each sequence and go to the end.
- Cannot detect domains

Smith-Waterman Algorithm

- Optimal Local Alignment
- Guaranteed to find all significant matches to a given query
- Takes the query sequence versus every sequence in the database
- Can be used with arbitrary scoring systems
- **COMPUTATIONALLY EXPENSIVE!!!**

Scoring Matrices

- Relatively simple for DNA-gap penalties or mismatches-can be made to look at Pu/Py
- Protein matches look also at similarity (leu/ileu)

Searching for Similarity

- FASTA
 - Search against user-defined search sets, DB or subsections
 - Only the single most similar region is shown

Nucleotide Uncertainties

Code Meaning (Base)

- A adenosine (A)
- C cytidine (C)
- G guanine (G)
- T thymidine (T)
- U uridine (U)
- R purine (G or A)
- Y pyrimidine (T or C)
- K keto (G or T)
- - gap(s) (none)

Code Meaning (Base)

- M amino (A or C)
- S strong (G or C)
- W weak (A or T)
- B not A (G or T or C)
- D not C (G or A or T)
- H not G (A or C or T)
- V not T (G or C or A)
- N any base (A or G or C or T)

Protein Scoring Matrices

- Chemical similarity: 210 pairs of aa
- Nearness in Genetic Code
- Chemical similarity, e.g.,
hydrophobicity
- Observed Substitution Schemes

AA Substitution Matrices

Rationale:

Certain amino acid substitutions commonly occur in related proteins (sometimes from different species). These provide the basis for amino acid substitution matrices, essentially a symbol comparison table.

More on Matrices

- A substitution matrix specifies a set of scores s_{ij} for replacing amino acid i by amino acid j .
- PAM: Percent Accepted Mutations
- BLOSUM Blocks Amino Acid Substitution Matrices

Amino Acid Symbols

A	Ala	alanine	R	Arg	Argine
B	Asx	Aspartic or asparagine	S	Ser	Serine
C	Cys	Cysteine	T	Thr	Threonine
D	Asp	Aspartic acid	U	Sec	Selenocysteine
E	Glu	Glutamic acid	V	Val	Valine
F	Phe	Phenylalanine	W	Trp	Tryptophan
G	Gly	Glycine	X	Xaa	Unknown or other aa
H	His	Histidine	Y	Tyr	Tyrosine
I	Ile	Isoleucine	Z	Glx	Glutamic or glutamine
K	Lys	Lysine			
L	Leu	Leucine			
M	Met	Methionine			
N	Asn	Asparagine			
P	Pro	Proline			
Q	Gln	Glutamine			

Observed AA Substitution Matrices

👉 PAM

👉 BLOSUM

PAM

- Log Odds scores are used
- The score of each pair $s(a,b)$ is defined as the log of the likelihood ratio of the transition probability M_{ab} (Mutation) versus the probability of a random occurrence of the amino acid b in the second sequence.

$$s(a,b) = \log M_{ab} / P_b$$

PAM: Point Accepted Mutation

- DAYHOFF et al.
- Observed residue replacement in related proteins
- GLOBAL alignment, closely related
- A model of molecular evolution
 - 1 PAM = average change in 1% of all amino acid possibilities (1% divergence)
- Other PAM matrices extrapolated from PAM1.

PAM continued

- ☞ TIME is NOT correlated with PAM
- ☞ Number of the matrix refers to evolutionary distance

Means different families of proteins evolve at different rates

PAM250

Table 2: The PAM250 matrix – an example of a matrix derived from observed substitutions

	A	B	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y	Z
A	2	0	-2	0	0	-4	1	-1	-1	-2	-1	0	1	0	-2	1	1	0	-6	-3	0	
B	0	2	-4	3	2	-5	0	1	-2	1	-3	-2	2	-1	1	-1	0	0	-2	-5	-3	2
C	-2	-4	12	-5	-5	-4	-3	-3	-2	-5	-6	-5	-4	-3	-5	-4	0	-2	-2	-8	0	-5
D	0	3	-5	4	3	-6	1	1	-2	0	-4	-3	2	-1	2	-1	0	0	-2	-7	-4	3
E	0	2	-5	3	4	-5	0	1	-2	0	-3	-2	1	-1	2	-1	0	0	-2	-7	-4	3
F	-4	-5	-4	-6	-5	9	-5	-2	1	-5	2	0	-4	-5	-5	-4	-3	-3	-1	0	7	-5
G	1	0	-3	1	0	-5	5	-2	-3	-2	-4	-3	0	-1	-1	-3	1	0	-1	-7	-5	-1
H	-1	1	-3	1	1	-2	-2	6	-2	0	-2	-2	2	0	3	2	-1	-1	-2	-3	0	2
I	-1	-2	-2	-2	1	-3	-2	5	-2	2	2	-2	-2	-2	-2	-1	0	4	-5	-1	-2	
K	-1	1	-5	0	0	-5	-2	0	-2	5	-3	0	1	-1	1	3	0	0	-2	-3	-4	0
L	-2	-3	-6	-4	-3	2	-4	-2	2	-3	6	4	-3	-3	-2	-3	-3	-2	2	-2	-1	-3
M	-1	-2	-5	-3	-2	0	-3	-2	2	0	4	6	-2	-2	-1	0	-2	-1	2	-4	-2	-2
N	0	2	-4	2	1	-4	0	2	-2	1	-3	-2	2	-1	1	0	1	0	-2	-4	-2	1
P	1	-1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6	0	0	1	0	-1	-6	-5	0
Q	0	1	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4	1	-1	-1	-2	-5	-4	3
R	-2	-1	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	1	6	0	-1	-2	2	-4	0	
S	1	0	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2	1	-1	-2	-3	0
T	1	0	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3	0	-5	-3	-1
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4	-6	-2	-2	
W	-6	-5	-8	-7	-7	0	-7	-3	5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	0	-6
Y	-3	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10	-4
Z	0	2	-5	3	3	-5	-1	2	-2	0	-3	-2	1	0	3	0	0	-1	-2	-6	-4	3

The non-standard amino acid B is used to refer to (D or N); 7 refers to (E or Q).

BLOSUM

- Block Substitution Matrix
- Henikoff and Henikoff, PNAS, 1992
- Number following indicates per cent identity within set, BLOSUM62=62% id
- Finds short, highly similar sequences (no gaps)

BLOSUM

- Matrices are directly calculated, based on observed alignments
- Greater numbers are lesser distances
- Usually best for local similarity searches
- BLOSUM62= DEFAULT FOR BLAST. If a distant relative, think about another matrix.

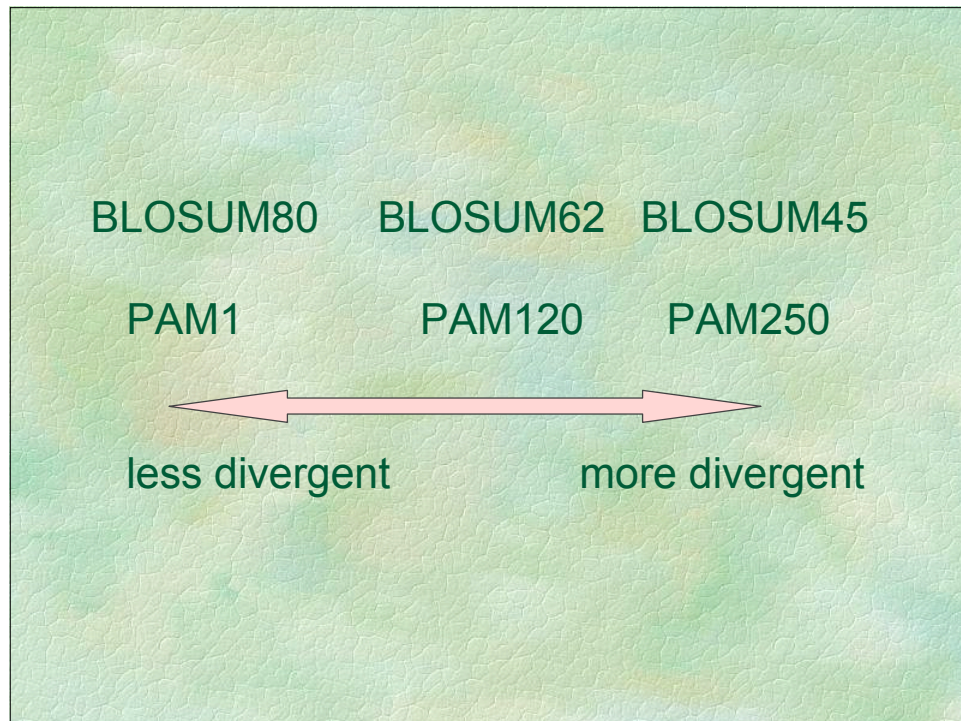
BLOSUM SCORING RULES

- Zero score means the frequencies of the pair in the database is that expected by chance
- A positive score means more frequent than chance
- Negative score means the pair is found less frequently than chance.

Blosum62

	A	C	D	E	F	G	H	
A	4	0	-2	-1	-2	0	-2	
C	0	9	-3	-4	-2	-3	-3	
D	-2	-3	6	2	-3	-1	-1	
E	-1	-4	2	5	-3	-2	0	
F	-2	-2	-3	-3	6	-3	-	
G	0	-3	-1	-2	-3			
H	-2	-3	-1	0				

BLOSUM 62



BLAST - Basic Local Alignment Sequence Tool

- Objective: find all local regions of similarity distinguishable from random
- Only local alignments permitted,
- Gaps permitted in version 2
- Statistically sound (Karlin and Altschul), but no guarantee of optimality

BLAST: Three Step Algorithm

- Compile a list of high scoring words of length w ($w=4$ for proteins, 12 for nucleic acids)
- Scan for word hits of score greater than threshold, T
- Extend word hit in both directions to find High Scoring Pairs with scores greater than S

Other BLAST Programs

- BLASTN: nucleic acid query to NA database
- BLASTP: Protein query to Protein database
- BLASTX: Translated nucleic acid query to Protein database
- TBLASTN: Protein query against (translated) nucleic acid database
- TBLASTX: Translated nucleic acid against translated nucleic acid database

OTHER BLAST VARIATIONS

- PSI-BLAST- Position Specific Iterated BLAST-use gapped BLAST, generate a Profile from multiple iterations used instead of the input and Distance Matrix
- MEGABLAST specifically designed to efficiently find long alignments between very similar sequences; the best tool to use to find the identical match to your query sequence.

Limitations to BLAST

- Needs islands of strong homology
- Limits on the combination of scoring and penalty values
- The variants (blastx, tblastn, tblastx) use 6-frame translation-miss sequences with frameshifts)
- Finds and reports ONLY local alignments


A WALK THROUGH BLAST

NCBI home

The screenshot shows the NCBI homepage with a blue header and a white background. The header includes the NCBI logo and navigation links: PubMed, Entrez, BLAST, OMIM, Books, TaxBrowser, and Structure. Below the header is a search bar with a dropdown menu for 'Nucleotide' and a 'Go' button. The main content area is divided into several sections:

- What does NCBI do?**: A paragraph explaining NCBI's mission, established in 1988, to create public databases, conduct research in computational biology, and develop software for analyzing genome data.
- PubMed Central**: A section highlighting free full-text articles from over 100 journals, linked to PubMed and fully searchable.
- Rat Genome Resource**: A section announcing the rat genome sequencing consortium, now available in NCBI's Map Viewer and BLAST.
- Mouse Genome**: A section featuring a map of the mouse genome and a link to the Mouse Genome browser.
- NCBI Newsletter**: A section encouraging users to subscribe to the newsletter for tips on improving BLAST results.
- Hot Spots**: A list of links to various resources, including Cancer genome anatomy project, Clusters of orthologous groups, Coffee Break, Electronic PCR, Gene expression omnibus, Genes and disease, Human genome resources, Human/mouse homology maps, LocusLink, Mammal genetics & genomics, Map Viewer, Mouse genome resources, NCBI Handbook, ORF Finder, Reference sequence project, Rattus resources, Serial analysis of gene expression, SKY/CGH database, SNP, Trace archive, UniGene, and VecScreen.

At the bottom of the page, there is a footer with the text 'Revised March 13, 2003' and a status bar showing 'Document: none'.



formatting **BLAST**

Nucleotide
Protein
Translations
Retrieve results for an RID

Your request has been successfully submitted and put into the Blast Queue.

Query = (8027 letters)

Your search was limited by an Entrez query: >(none)

The request ID is


or

The results are estimated to be ready in 21 minutes but may be done sooner.

Please press "FORMAT!" when you wish to check your results. You may change the formatting options for your request via the form below and press "FORMAT!" again. You may also request results of a different search by entering any other valid request ID to see other recent jobs.

Back
Forward
Stop
Refresh
Home
Search
Favorites
Media
History
Mail
Print
Edit
Discuss
Go
Links

Address
http://www.ncbi.nlm.nih.gov/blast/Blast.cgi



results of **BLAST**

BLASTN 2.2.6 [Apr-09-2003]

Reference:
Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

RID: 1052261999-028139-8384

Query= gi|16933541|ref|NM_002026.1| Homo sapiens fibronectin 1 (FN1), transcript variant 1, mRNA (8027 letters)

Database: All GenBank+EMBL+DDBJ+PDB sequences (but no EST, STS, GSS, or phase 0, 1 or 2 HTGS sequences)
1,747,911 sequences; 8,529,560,197 total letters

If you have any problems or questions with the results of this search please refer to the [BLAST FAQs](#)

[Taxonomy reports](#)

Distribution of 286 Blast Hits on the Query Sequence

Mouse-over to show define and scores. Click to show alignments

Color Key for Alignment Scores

<40

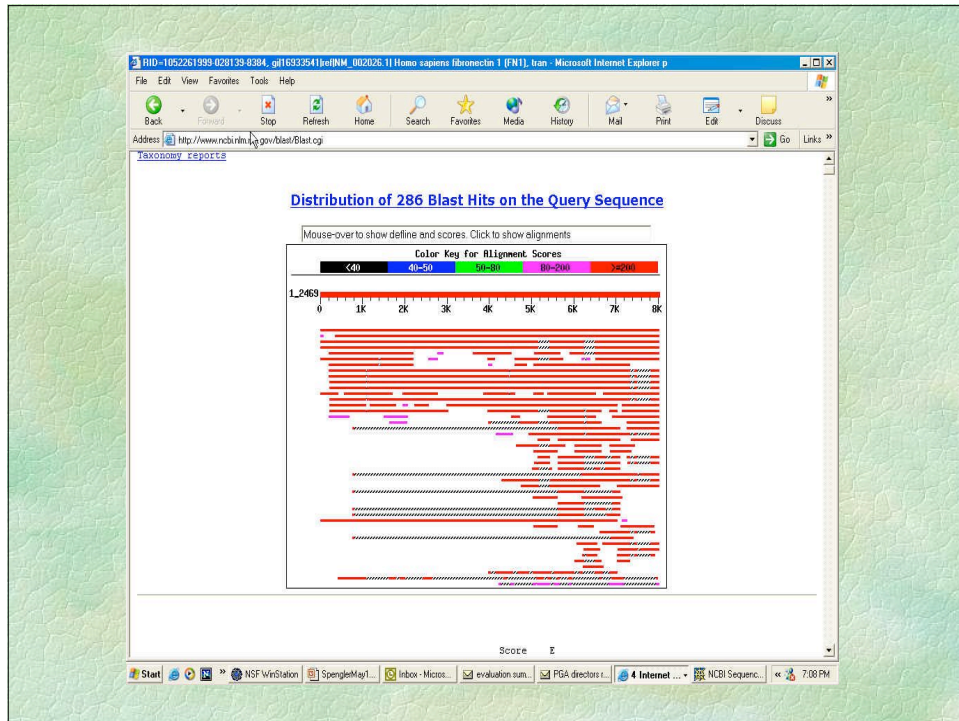
40-50

50-80

80-200

>200

Start
NSF WinStation
SpenglerMay1...
Inbox - Micros...
evaluation sum...
PGA directors f...
4 Internet ...
NCBI Sequenc...
7:07 PM



2 RID-1052261999-020139-8384_gi16933541|ref|NM_002026.1| Homo sapiens fibronectin 1 (FN1), tran - Microsoft Internet Explorer p

File Edit View Favorites Tools Help

Address <http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>

Sequences producing significant alignments:

Accession	Description	Score	E (bits)	Value
gi16933541 ref NM_002026.1 	Homo sapiens fibronectin 1 (FN1), tran	1.591e+04	0.0	L U
gi131961 emb X02761.1 HSFIB1	Human mRNA for fibronectin (FN1)	1.517e+04	0.0	L U
gi21733353 emb AL832771.1 HSM804082	Homo sapiens mRNA; cDN...	1.014e+04	0.0	L U
gi21732747 emb AL832202.1 HSM803509	Homo sapiens mRNA; cDN...	1.014e+04	0.0	L U
gi27227742 emb J0535086.1 HSA535086	Homo sapiens mRNA for ...	3943	0.0	L
gi11493493 gb AF130095.1 AF130095	Homo sapiens clone FLC05...	3003	0.0	
gi16877226 gb BC016875.1 	Homo sapiens, clone IMAGE386658...	2970	0.0	
gi1826961 gb H10905.1 HUMFNC	Human cellular fibronectin mRNA	2904	0.0	L U
gi16933543 ref NM_054034.1 	Homo sapiens fibronectin 1 (FN1)	2730	0.0	L U
gi19506702 ref NM_019143.1 	Rattus norvegicus Fibronectin 1...	2718	0.0	L U
gi56163 emb X15906.1 FNFIBRON	Rat mRNA for fibronectin	2718	0.0	L U
gi28479105 ref XM_129845.3 	Mus musculus fibronectin 1 (Fn...	2668	0.0	L
gi26105748 dbj AK090135.1 	Mus musculus 5 months female bo...	2652	0.0	L
gi13543399 gb BC005858.1 BC005858	Homo sapiens, clone MGC...	2627	0.0	L
gi26105744 dbj AK090130.1 	Mus musculus 5 months female bo...	2565	0.0	L
gi1124966 emb X93167.1 MMFIB1	M.musculus mRNA for fibronectin	2440	0.0	L
gi4096861 gb U42594.1 HSU42594	Human fibronectin (FN1) mRN...	2371	0.0	L
gi12053816 emb J0276395.1 HSA276395	Homo sapiens mRNA for ...	2296	0.0	L U
gi4096845 gb U42404.1 HSU42404	Human fibronectin (FN1) mRN...	2234	0.0	L
gi4204942 gb U60067.1 HSU60067	Human fibronectin mRNA, par...	2228	0.0	
gi4096851 gb U42457.1 HSU42457	Human fibronectin (FN1) mRN...	1820	0.0	L
gi29835173 gb BC051082.1 	Mus musculus, Similar to Fibrone...	1790	0.0	
gi10439658 dbj AK026737.1 	Homo sapiens cDNA: FLJ23084 fis...	1733	0.0	L U
gi21753154 dbj AK094153.1 	Homo sapiens cDNA FLJ36834 fis...	1731	0.0	
gi4204944 gb U60068.1 HSU60068	Human fibronectin mRNA, alt...	1669	0.0	
gi40968591 gb U42593.1 HSU42593	Human fibronectin (FN1) mRN...	1612	0.0	L

[RID-1052261989-028139-9384](#), [gi16933541|refNM_002026.1|Homo sapiens fibronectin 1 \(FN1\), tran](#) - Microsoft Internet Explorer p

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss Go Links

Address <http://www.ncbi.nlm.nih.gov/blast/blast.cgi#26105748>

>[gi126105748|db11AK090135.1](#) Mus musculus 5 months female bone marrow stroma cell CRL-2028 SR-4987
 cDNA, RIKEN full-length enriched library,
 clone:G431004B19 product:fibronectin 1, full insert
 sequence
 Length = 835

Score = 2652 bits (1338), Expect = 0.0
 Identities = 2544/2944 (86%), Gaps = 4/2944 (0%)
 Strand = Plus / Plus

Query: 1110 cggatctggcccttcaccgatgtgtgcagctgtttaccacccgagcctcaccacca 1169
 ||||| ||||| ||||| || ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
 Sbjct: 1036 cggatctggctcttcactgatgtccgaacagctattaccacccgagactcaccacca 1095

Query: 1170 gccctccctccatggccactgtgtcacagacagtggtgtggtctactctgtgggagatgca 1229
 ||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
 Sbjct: 1096 gccctccctccatggccactgtgtcacagacagtggtgtggtctactctgtgggagatgca 1155

Query: 1230 gtggctgaagacacacaggaataagcaaatgctttgcacgtgctgggcaacggagtcag 1289
 ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
 Sbjct: 1156 gtggctgaagtcgcaaggaacacagcaaatgctgtgcacgtgctgggcaacggagtcag 1215

Query: 1290 ctgccaagagacagctgttaaccacagcttaccgttgcaactcaaatggagagccatgtgt 1349
 ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
 Sbjct: 1216 ctgccaagagacagcgtgacccagacttatgttgcaattcaaacgggagcgcctgtgt 1275

Query: 1350 cttaccattacctaacaatggcaggacgttctactcctgcaccacgaaggcgacagga 1409
 || || ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
 Sbjct: 1276 cctcccgctcactacaacggtaggacctctctactcctgcaccacgaaggcgaggaaga 1335

Query: 1410 cggacatctttgtgagcagacacactcgcaattatgagcaggaccagaaatactctttctg 1469
 ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||

Back Forward Reload Home Search Netscape Print Security Shop Stop

Bookmarks Location <http://www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?val=16933541&db=Nucleotide&dopt=GenBank> What's Related

Members WebMail Connections BioJournal SmartUpdate Mixplace Evolution Home Raising Duncan DCA - Ronald Re Physical Therap M

NCBI Nucleotide Protein Genome Structure PMC Taxonomy OMIM Books

Search Nucleotide for Go Clear

Limits Preview/Index History Clipboard Details

Display default Show: 1 Send to File Get Subsequence

1: NM_002026. Homo sapiens fibr...[gi:16933541]

LOCUS FN1 8027 bp mRNA linear PRI 03-APR-2003
 DEFINITION Homo sapiens fibronectin 1 (FN1), transcript variant 1, mRNA.
 ACCESSION NM_002026
 VERSION NM_002026.1 GI:16933541
 KEYWORDS .
 SOURCE Homo sapiens (human)
 ORGANISM [Homo sapiens](#)
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
 Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.
 REFERENCE 1 (bases 1 to 8027)
 AUTHORS Rameshwar,P., Oh,H.S., Yook,C., Gascon,P. and Chang,V.T.
 TITLE Substance p-fibronectin-cytokine interactions in myeloproliferative disorders with bone marrow fibrosis
 JOURNAL Acta Haematol. 109 (1), 1-10 (2003)
 MEDLINE 22372683
 PUBMED 12486316
 REMARK GeneRIF: REVIEW: Substance p-fibronectin-cytokine interactions in myeloproliferative disorders with bone marrow fibrosis.
 REFERENCE 2 (bases 1 to 8027)
 AUTHORS Arici,M., Brown,J., Williams,M., Harris,K.P., Walls,J. and Brumekitt,N.J.

Document Done

Transcript Variant: This variant (1) lacks an exon containing the EDII (extra domain II) region. It encodes a 2355 aa isoform of the protein.

COMPLETENESS: full length.

FEATURES

source

1..8027

/organism="Homo sapiens"

/mol_type="mRNA"

/db_xref="taxon:9606"

/chromosome="2"

/map="2q34"

gene

1..8027

/gene="FN1"

/note="synonyms: FN, CIG, FINC, LETS"

/db_xref="LocusID:2335"

/db_xref="MIM:135600"

CDS

268..7335

/gene="FN1"

/note="isoform 1 is encoded by transcript variant 1; cold-insoluble globulin"

/codon_start=1

/product="fibronectin 1 isoform 1 preproprotein"

/protein_id="NP_002017.1"

/db_xref="GI:16933542"

/db_xref="LocusID:2335"

/db_xref="MIM:135600"

/translation="MLRGPGQLLLAVQCLGTAVPSTGASKSKRQAQQMVQPSVVA
VSQKPGCYDNGKHQYINQWERTYLGNALVCTCYGSGRGFNCSEKPEAETCFDKYT
GNTYKVDITYERPKDSIMWCTCIGAGRGISCTIANRCHGGQSYKIGDTWRPPHET
GGYMLCVCVLGNGKGWTCRPLAEKCFDHAAGTSYVGETWEKPYQGMVMDCTCLGE
GSGRLTCTSRNRCNDQDTRTSYRIGDTWSEKDNKGNLLQICTGNGRGWKECHRTSV
QTTSGSGSPFTDVRAAVTQPHQPPYKGVCTDSGVVSVGMQLKTQGNKQMLCT
CLGNGVSCQETAVTQTYGNSNGEPCLVPTFYNGRTFYSCCTEGRQDGLWCSTTSNY
EQDKYSFCTDHTLVQTRGNSNGALCHFFLYNNHNYDCTSEGRDMKWCQTQ

with fibrin-, heparin-, s.aureus-binding activity"

misc feature

1189..1299

/gene="FN1"

/note="region of internal homology I (1 subunit); domain with collagen-binding activity"

misc feature

1300..1674

/gene="FN1"

/note="region of internal homology II (2 subunits) domain with collagen-binding activity"

misc feature

1675..2091

/gene="FN1"

/note="region of internal homology I (3 subunits) domain with collagen-binding activity"

misc feature

2092..2367

/gene="FN1"

/note="region of internal homology III (1 subunit)"

misc feature

2422..6240

/gene="FN1"

/note="region of internal homology III (14 subunits)"

misc feature

2694..3522

/gene="FN1"

/note="domain with DNA-binding activity"

misc feature

4063^4064

/gene="FN1"

/note="alternatively spliced exon: EDII region; see Accession M18178.1"

misc feature

4837..4848

/gene="FN1"

/note="cell binding site"

misc feature

5158..5427

/gene="FN1"

/note="alternatively spliced exon: EDI region"

misc feature

5428..6240

/gene="FN1"

Back Forward Reload Home Search Netscape Print Security Shop

Bookmarks Location: <http://www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?val=16933541&db=Nucleotide&dopt=GenBank> What's Related

Members WebMail Connections BioJournal SmartUpdate Mikiplace Evolution: Home Raising Duncan DCA - Ronald Re Physical Therap M

BASE COUNT 2130 a 2111 c 1973 g 1813 t

ORIGIN

```

1  acgcccgcgc  cggctgtgct  gcacaggggg  aggagaggga  accccaggcg  cgagcgggaa
61  gaggggacct  gcagccaaaa  cttctctggt  cctctgcato  ccttctgtcc  ctcaccocgt
121  ccccttcccc  accctctggc  ccccaccttc  ttggaggcga  caaccccggg  gaggcattag
181  aagggatttt  tcccgcagtt  gcgaaggaaa  gcaaaccttg  tggcaacttg  cctcccggtg
241  cgggcgtctc  tccccaccgc  tctcaacatg  cttagggggc  cggggcccg  gctcgtctg
301  ctggccctcc  agtgcctggg  gacagcgggt  cctccaccgg  gagcctcgaa  gagcaagagg
361  caggctcagc  aaatggttca  gccccagtc  ccggtggctg  tcagtcaaag  caagcccggt
421  tgttatgaca  atggaaaaaa  ctatcagata  aatcaacagt  gggagcggac  ctacctaggc
481  aatgcgttgg  ttgttaactg  ttatggagga  agccgaggtt  ttaactcgga  gaggtaaacct
541  gaagctgaag  agacttgctt  tgacaagtag  actgggaaca  cttaccgagt  ggtgacacct
601  tatgagcgtc  ctaaaagact  catgatctgg  gactgtacct  gcactggggc  tgggcgaggg
661  agaataagct  gtaccatcgc  aaaccgctgc  catgaagggg  gtcagtccta  caagatttgt
721  gacacctgga  ggagaccaca  tgagactggt  ggttacatgt  tagagtgtgt  gtgcttttgt
781  aatggaaaag  gagaatggac  ctgcaagccc  atagctgaga  agtgttttga  tcatgcttgt
841  gggactctct  atgctgtcgg  agaaacctgg  gagaagccct  accaaggctg  gatgatggta
901  gattgtactt  gctcgggaga  aggcagcgga  cgcatacttt  gcacttctag  aaatagatgc
961  aacgatcagg  acacaaggac  atcctataga  attggagaca  cctggagcaa  gaagataaat
1021  cgaggaaaac  tgcctccagt  catctgcaca  ggcaacggcc  gaggagagtg  gaagtgtgag
1081  aggcacacct  ctgtgcagac  cacatcgagc  ggcattggcc  ccttcaccga  tgttcgtgca
1141  gctgtttacc  aaccgcagcc  tcacccccag  cctctcctct  atggccactg  tgtcacagac
1201  agtgggttgg  tctactctgt  ggggatgcag  tggctgaaga  cacaaggaaa  taagcaaatg
1261  ctttgacagt  gctcgggcaa  cggagtcagc  tgccaagaga  cagctgtaac  ccagacttac
1321  ggtggcaact  caaatggaga  gccatgtgtc  ttaccattca  cctacaatgg  caggacgttc
1381  tactcctgca  ccacggaagg  gcgacaggac  ggacattctt  ggtgcagcac  aacttcgaat
1441  tatgagcagg  accgaaata  cttcttctgc  acagaccaca  ctgttttgtt  tcagactcga
1501  ggaggaaatt  ccaatggtgc  cttgtgccac  ttccctctcc  tatacaacaa  ccacaattac
1561  actgattgca  cttctgaggg  cagaagagac  aacatgaagt  ggtgtgggac  cacacagaac
1621  tatgatgcgc  accgaaagtt  tgggtctctc  cccatggctg  cccacgagga  aatctgcaca
1681  accaatgaag  ggtcactgta  ccgcatctga  gatcagtggt  ataagcagca  tgacatgggt
1741  cacatgatga  ggtgcacgtg  tgttgggaat  ggtcgtgggg  aatggacatg  cattgcctac
1801  tcgcagcttc  gagatcagtg  cattgttgat  gacatcactt  acaatgtgaa  cgacacattc
1861  cacaagcgtc  atgaagaggg  gcacatgctg  aactgtacat  gcttcgggta  gggtcggggc
1921  cccacacgct  cccacacgct  cccacacgct  cccacacgct  cccacacgct  cccacacgct

```

Document: Done

What happens with changes?

- Word size
- Substitution Matrix
- Expected Value
- Using protein or translated sequence rather than nucleotide

Options for BLAST

Options for advanced blasting

Submit by [entrez query](#) or select from:

[Composition-based statistics](#) ☐

[Choose filter](#) ☐ Low complexity ☐ Mask for lookup table only ☐ Mask lower case

[Expect](#)

[Word Size](#)

[Matrix](#) Gap Costs

[PSSM](#)

[Other advanced](#)

[PHI pattern](#)

The Word-Size Parameter

A word is any short sequence less than or equal to

- Protein 1-2
 - Nucleotide 1-12
- Higher word Size
- Faster
 - Less Sensitive
 - More Selective

Expected Value

- Number of hits with a score better than “S” expected by chance in a DB of a given size
- Default value of threshold is 10.
- For a short sequence, increase the E value (say to 1000)

BLAST RULES OF THUMB

- For short amino acid sequences (20-40), 50% identity happens by chance
- If A and B are homologous, and B and C are homologous, then A and C are, even if you can't see it.
- You can get similarity in the absence of homology for low complexity, transmembrane and coiled-coil regions. These have to be eliminated by you, but you MAY want them.

TAKE-HOMES

- ☛ Use an up-to-date database; repeat often
- ☛ Choose a fast algorithm
- ☛ Use the most recent version
- ☛ Work at the protein level--for a small amount of evolutionary change, DNA sequence contains less information about homology
- ☛ Respect your own *intuition*

BLAST Significance

- ☛ If you change scoring systems, you can still compare search results if you normalize the score.

$S' = (\lambda S - \ln K) / \ln 2$. λ and K are associated with the scoring system.

S' , with a given E , is significant if it is greater than $\log N/E$, N the size of the search space.

Searching for Similarity

• FASTA

- Search against user-defined search sets, DB or subsections
- Only the single most similar region is shown

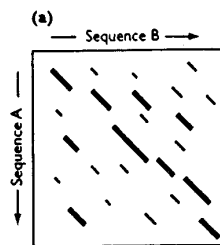
FASTA: WHY USE IT?

- Allow alignments to shift frames

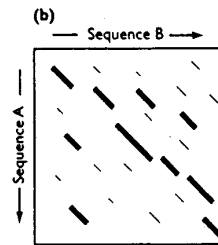
FASTA: FAST Alignment

- <http://alpha10.bioch.virginia.edu/fasta/>
- <http://www2.ebi.ac.uk/fasta3>
- <http://workbench.sdsc.edu>
- Rapid Global alignment
- Not a strong mathematical basis

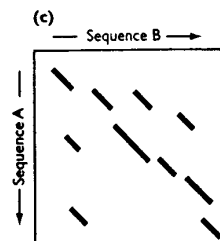
FASTA Algorithm



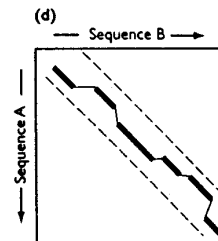
Find runs of identities.



Re-score using PAM matrix.
Keep top scoring segments.



Apply 'joining threshold' to eliminate segments that are unlikely to be part of the alignment that includes highest scoring segment.



Use dynamic programming to optimise the alignment in a narrow band that encompasses the top scoring elements.

LALIGN

- Essentially a FASTA derivative for local alignments
- Compares two proteins to identify regions of similarity
- Will report several sequence alignments within a given sequence
- Works for internal repeats that are missed by FASTA because of gaps.

SITEs for LALIGN

- <http://fasta.bioch.virginia.edu/fasta/lalign.htm>
- <http://xylian.igh.cnrs.fr/bin/lalign-guess.cgi>
- <http://biowb.sdsc.edu> (registration necessary but painless)
- PALIGN
<http://fasta.bioch.virginia.edu/fasta/palign.htm>
(plots a graph of the areas of alignment)

Gene Ontologies GO

A gene ontology is a controlled vocabulary used to describe the biology of a gene product in any organism, designed to allow both attribution and querying at different levels of granularity, facilitating queries across participating databases.

A step toward unifying biological databases but not sufficient.

<http://www.geneontology.org>

Components of GO

A gene product is a physical thing (protein, RNA, can have small molecules associated to make a gene product group.

Attributes of Gene Products

- 🐼 **Molecular Function**-what something does
- 🐼 **Biological process**-a biological objective, like growth or pyrimidine metabolism
- 🐼 **Cellular Component**-part of a cell, ER, nucleus etc.

Ontology Representations

- A network, a directed acyclic graph (DAG), in which terms have multiple parents and multiple relationships to parents.
- Relationships connecting terms include is-a, part-of, Yeast, Fly, Mouse, Arabidopsis, Worm,

EVIDENCE CODES

- IC Inferred by Curator
- IDA Inferred by Direct Assay
- IEA Inferred by Electronic Annotation
- IEP Inferred from expression pattern
- IGI Inferred from genetic interaction
- IMP Inferred from mutant phenotype
- IPI Inferred from physical interaction
- ISS Inferred from sequence or structure similarity
- NAS Non-traceable author statement
- ND No biological data available
- TAS Traceable author statement
- NR Not recorded

Evidence relationships

TAS/IDA

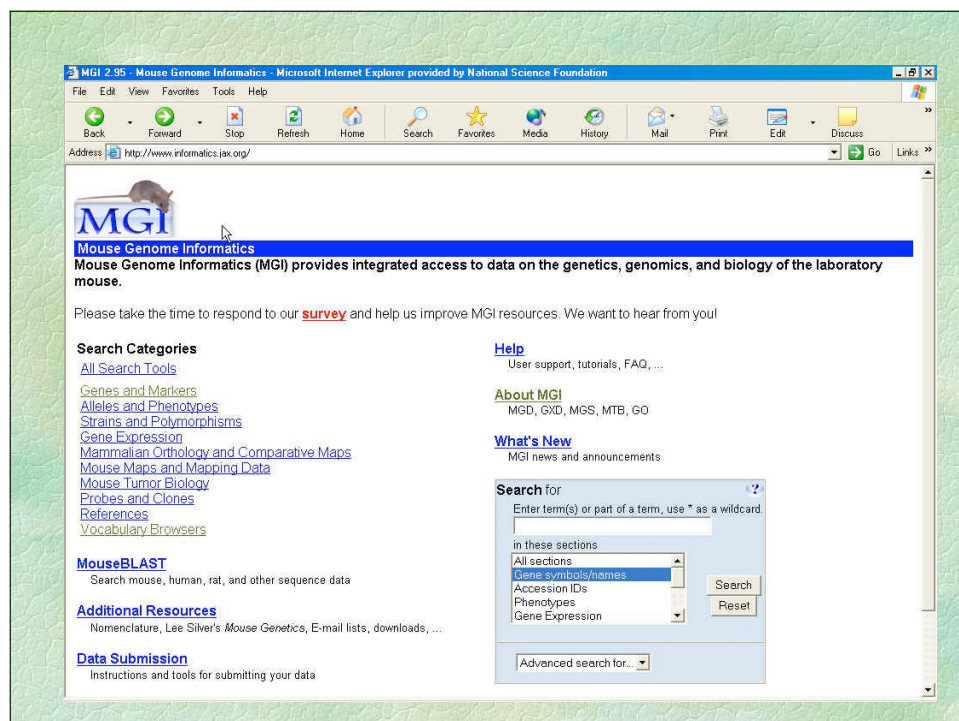
IMP/IGI/IPI


ISS/IEP

NAS

IEA

Not a rigid hierarchy.





[MGI Home](#) [Help](#)

Search for: Go

in these sections

[All sections](#)
[Gene symbols/names](#)
[Accession IDs](#)
[Phenotypes](#)
[Gene Expression](#)

Advanced search for...


Search Categories (All tools)

[Genes/Markers](#)
[Alleles/Phenotypes](#)
[Strains/Chromosomes](#)
[Expression](#)
[Comparative Maps/Data](#)
[Mouse Maps/Genes](#)
[Mouse Tumor Biology](#)
[Protein/Clones](#)
[References](#)
[Vocabulary Browser](#)
[Gene Ontology \(GO\)](#)
[Anatomical Dictionary](#)
[Phenotype Classifications](#)

[MouseBLAST](#)

Additional Resources

[China Thesaurus](#)
[Funding Information](#)
[Warranty Disclaimer](#)
[Copyright Notice](#)
[Send questions and comments to User Support](#)


 last database update
 02/20/2003
 MGI 2.9

Gene Ontology Browser Term Detail

GO term: **blood coagulation factor VII**
 GO id: **GO:0003802**
 Definition: **Catalysis of the selective cleavage of one Arg-Ile bond in factor X to form factor Xa.**
 Number of paths to term: 3

ⓘ denotes an 'is-a' relationship
 ⓘ denotes a 'part-of' relationship

Gene_Ontology

- ⓘ [molecular_function](#)
- ⓘ [defense/immunity protein](#)
- ⓘ [blood coagulation factor](#)
 - ⓘ [blood coagulation factor IX](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)
 - ⓘ [blood coagulation factor X](#)
 - ⓘ [blood coagulation factor XII](#)
 - ⓘ [blood coagulation factor XIII](#)
 - ⓘ [plasmin kallikrein](#)
 - ⓘ [protein C \(activated\)](#)
 - ⓘ [protein-cytoskeleton gamma-cytoplasmic transmembrane](#)
 - ⓘ [thrombin](#)

Gene_Ontology

- ⓘ [molecular_function](#)
- ⓘ [defense/immunity protein](#)
- ⓘ [complement activity](#)
 - ⓘ [alternative-complement pathway C3/C5 convertase](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)
 - ⓘ [classical-complement pathway C3/C5 convertase](#)
 - ⓘ [complement component C1x](#)
 - ⓘ [complement component C1s](#)
 - ⓘ [complement factor D](#)
 - ⓘ [complement factor H](#)
 - ⓘ [complement factor I](#)

Gene_Ontology

- ⓘ [molecular_function](#)
- ⓘ [enzyme](#)
- ⓘ [hydrolase](#)
- ⓘ [peptidase](#)
- ⓘ [serine-type peptidase](#)
- ⓘ [serine-type endopeptidase](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)

Gene Ontology Browser Term Detail

GO term: **blood coagulation factor VII**
 GO id: **GO:0003802**
 Definition: **Catalysis of the selective cleavage of one Arg-Ile bond in factor X to form factor Xa.**
 Number of paths to term: **3**

ⓘ denotes an 'is-a' relationship
 ⓘ denotes a 'part-of' relationship

Gene_Ontology

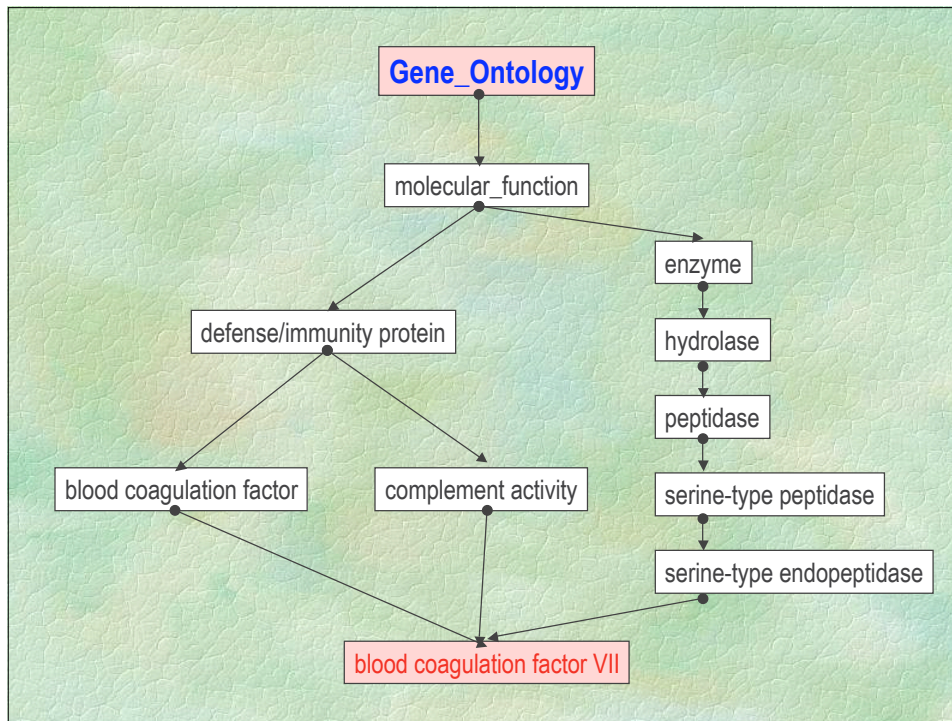
- ⓘ [molecular_function](#)
- ⓘ [defense/immunity protein](#)
- ⓘ [blood coagulation factor](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)

Gene_Ontology

- ⓘ [molecular_function](#)
- ⓘ [defense/immunity protein](#)
- ⓘ [complement activity](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)

Gene_Ontology

- ⓘ [molecular_function](#)
- ⓘ [enzyme](#)
- ⓘ [hydrolase](#)
- ⓘ [peptidase](#)
- ⓘ [serine-type peptidase](#)
- ⓘ [serine-type endopeptidase](#)
 - ⓘ [blood coagulation factor VII \[GO:0003802\] \(*L.genes_1 annotations*\)](#)



GO Browsers	
AmiGO from BDGP	<ul style="list-style-type: none"> With AmiGO, you can search for a GO term and view all gene products annotated to it, or search for a gene product and view all its associations. You can also browse the ontologies to view relationships between terms as well as the number of gene products annotated to a given term. AmiGO accesses the GO MySQL database (see below); the browser and documentation are available from http://www.godatabase.org/dev/
MGI GO Browser	<ul style="list-style-type: none"> With the MGI GO Browser, you can search for a GO term and view all mouse genes annotated to the term or any subterms. You can also browse the ontologies to view relationships between terms, term definitions, as well as the number of mouse genes annotated to a given term and its subterms. The MGI GO browser directly accesses the GO in the MGI database where mouse gene annotations, are updated nightly. The version of the GO used is obtained nightly from the GO ftp site.
QuickGO at EBI	<ul style="list-style-type: none"> With QuickGO, a GO browser integrated into InterPro at the EBI, you can search for a GO term to see its relationships and definition, as well as any available mappings to SWISS-PROT keywords, to the Enzyme Classification or Transport Classification databases, or to InterPro entries. Use documentation is available from the manual and the FAQ.
EP GO Browser	<ul style="list-style-type: none"> The EP-GO browser is built into EBI's Expression Profiler, a set of tools for clustering, analysis and visualization of gene expression and other genomic data. With it, you can search for GO terms and identify gene associations for a node, with or without associated subnodes, for the organism of your choice.
GoFish	<ul style="list-style-type: none"> The GoFish program, available as a Java applet, allows the user to construct arbitrary Boolean queries using GO attributes, and orders gene products according to the extent they satisfy such queries. GoFish also estimates, for each gene product, the probability that they satisfy the Boolean query. Developed by the Roth lab at Harvard.
GenNav	<ul style="list-style-type: none"> GenNav is a GO browser developed at NLM. It searches GO terms and annotated gene products, and provides a graphical display of a term's position in the GO DAG.
GeneOntology@RZPD	<ul style="list-style-type: none"> With the GeneOntology@RZPD tool at the Resource Center/Primary Database (RZPD) in Germany, you can search for GO identifiers associated with UniGene Clusterids, Genes (Name/Symbol) and Clones provided by the RZPD. You can also search for UniGene Clusters, Genes and Clones annotated with a certain GO identifier or a combination of GO identifiers. So far, GO annotations for human and mouse genes/clones are linked.
ProToGO	<ul style="list-style-type: none"> ProToGO, developed at the Hebrew University in Jerusalem, searches the GOA@EBI and CompuGen annotation datasets. The output is a graphical view of the relevant sub-graph of GO, containing those GO terms assigned to the query proteins. Documentation is provided.
CGAP GO Browser	<ul style="list-style-type: none"> With the GO browser at the The Cancer Genome Anatomy Project, you can browse through the GO vocabularies, and find human and mouse genes assigned to each term. The help documentation is at: http://cgap.nci.nih.gov/Genes/AboutGO.
DAG-Edit	
DAG-Edit	<p>This Java application provides an interface to browse, query and edit GO or any other vocabulary that has a DAG data structure. The most current version of DAG-Edit can be downloaded from the publicly accessible source repository at SourceForge. Help documentation to use the program can also be downloaded from this site (.pdf or .html formats) or is available here: http://www.geneontology.org/doc/dagedit_userguide/dagedit.html</p>
GO Database	
GO Database	<p>API documentation, schema diagrams and full descriptions of all tables for the MySQL database developed and maintained by BDGP, http://www.godatabase.org/dev/database/</p>

GO Database	
GO Database	API documentation, schema diagrams and full descriptions of all tables for the MySQL database developed and maintained by BDGP, http://www.godatabase.org/dev/database/
Other GO Tools	
GO Term Finder	<ul style="list-style-type: none"> The GO Term Finder at SGD searches for significant shared GO terms, or parents of the GO terms, used to annotate budding yeast gene products.
GO Term Mapper	<ul style="list-style-type: none"> The GO Term Mapper at SGD maps the specific, granular GO terms used to annotate a list of budding yeast gene products to corresponding GO Slim terms (i.e. more general parent GO terms; uses the SGD GO Slim set).
Manatee	<ul style="list-style-type: none"> Manatee is a web-based gene evaluation and genome annotation tool developed at TIGR. Manatee can store and view annotation for prokaryotic and eukaryotic genomes. The Manatee interface allows biologists to quickly identify genes and make high quality functional assignments, such as GO classifications, using search data, paralogous families, and annotation suggestions generated from automated analysis.
PubSearch	<ul style="list-style-type: none"> PubSearch is a web-based literature curation tool developed at TAIR and available via GMOD. It allows curators to search and annotate genes to keywords from articles. It has a simple, MySQL database backend and uses a set of Java Servlets and JSPs for querying, modifying, and adding gene, gene-annotation, and literature information. A demo is available.
SOURCE	<ul style="list-style-type: none"> SOURCE, developed by the Stanford Microarray Database (SMD) team, compiles information from several publicly accessible databases, including UniGene, dbEST, Swiss-Prot, GeneMap99, RHDb, GeneCards and LocusLink. GO terms associated with LocusLink entries appear in SOURCE.
MAPPFinder	<ul style="list-style-type: none"> MAPPFinder is an accessory program for GenMAPP. This program allows users to query any existing GenMAPP Expression Dataset Criterion against GO gene associations and GenMAPP MAPPs (microarray pathway profiles). The resulting analysis provides the user with results that can be viewed directly upon the Gene Ontology hierarchy and within GenMAPP, by selecting terms or MAPPs of interest.
FatiGO	<ul style="list-style-type: none"> FatiGO is a web interface for clustering DNA microarray data and simple datamining using GO. datamining consists of the assignment of the most characteristic Gene Ontology term to a cluster. GO terms are related to Unigene Human and Mouse Cluster Ids and Saccharomyces Genome Database.
Onto-Express	<ul style="list-style-type: none"> Onto-Express searches the public databases and returns tables that correlate expression profiles with the cytogenetic gene locations, the biochemical and molecular functions, the biological processes, cellular components and cellular roles of the translated proteins. (Registration required; free for academics.)
Genes2Diseases	<ul style="list-style-type: none"> Genes2Diseases is a database of candidate genes for mapped inherited human diseases, developed by the Bork group at the European Molecular Biology Laboratory (EMBL). The database is generated using an analysis of relations between phenotypic features and chemical objects, and from chemical objects to protein function (Gene Ontology) terms, based on the whole MEDLINE and RefSeq databases. Can be used to view all GO terms associated with a particular genetically inherited disease.

GO and Links

GO

Linkout to SGD

SGD Quick Search: Submit [Site Map](#) [Help](#) [Full Search](#) [Home](#) [Help](#)

[Community Info](#) [Submit Data](#) [BLAST](#) [Primers](#) [PatMatch](#) [Gene/Seq Resources](#) [Virtual Library](#) [Contact SGD](#)

RPE1/YJL121C

[Help](#)

Alternative single page format [new](#)

RPE1 BASIC INFORMATION

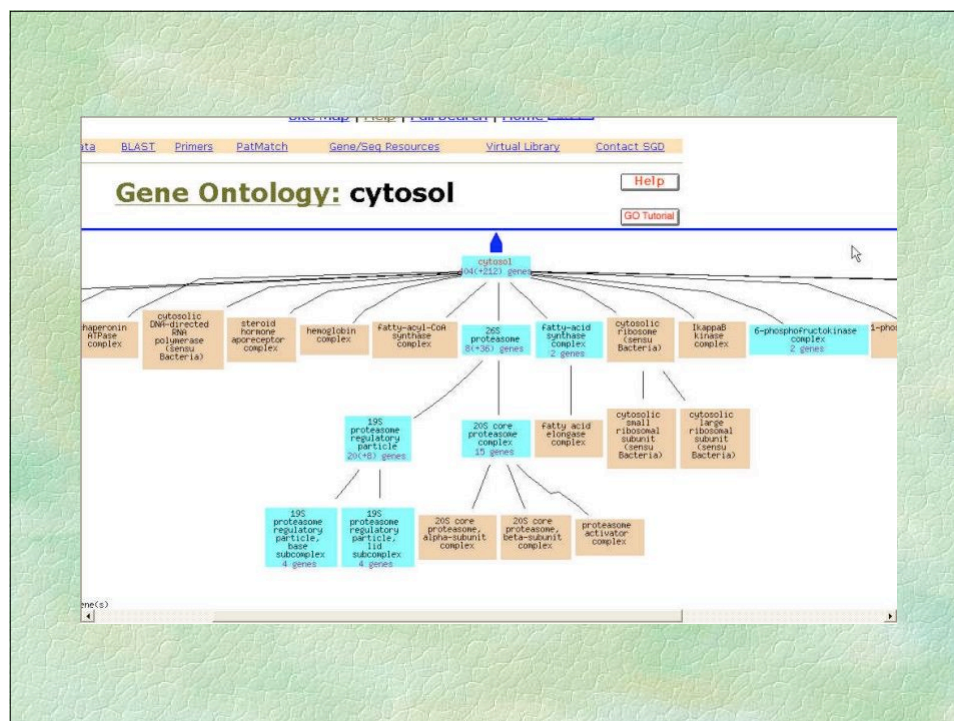
Standard Name	RPE1
Alias	EPT1, POS18
Systematic Name	YJL121C
Feature Type	ORF
GO Annotations	RPE1 GO evidence and references
Molecular Function	<ul style="list-style-type: none"> ribulose-phosphate 3-epimerase activity
Biological Process	<ul style="list-style-type: none"> pentose-phosphate shunt
Cellular Component	<ul style="list-style-type: none"> cytosol
Description	D-ribulose-5-phosphate 3-epimerase
Gene Product	D-ribulose-5-phosphate 3-epimerase
Phenotype	RPE1 Phenotype details and references
Systematic deletion	<ul style="list-style-type: none"> Order mutant strains used in the systematic deletion project
Free text	<ul style="list-style-type: none"> violate Exhibits growth defect on a non-fermentable (respiratory) carbon source. Null mutants are viable but show no ribulose-5-phosphate epimerase activity, cannot grow on D-xylulose, and are sensitive to hydrogen peroxide
Sequence Coordinates	ChrX: coordinates 190790 to 190074 [ORF Map]
Exon	1 - 717 190790 - 190074
External Links	MIPS YPP SwissProt Entrez Protein Entrez Neighbors EBI/EMBL PDB/US Entrez RefSeq Biozyme Kyoto
Primary SGDID	S0003657

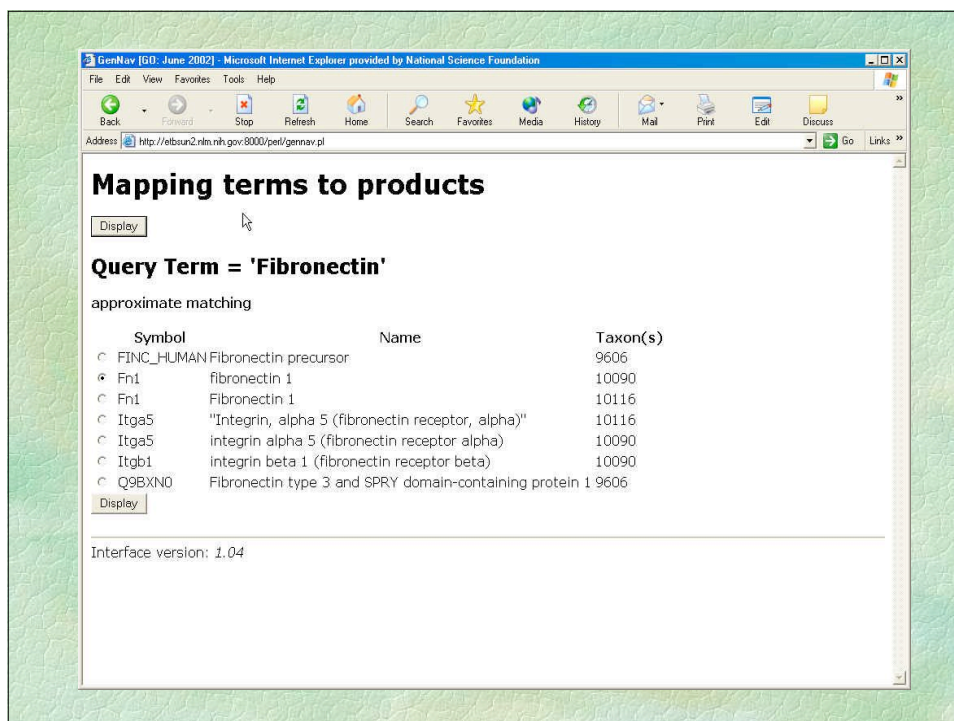
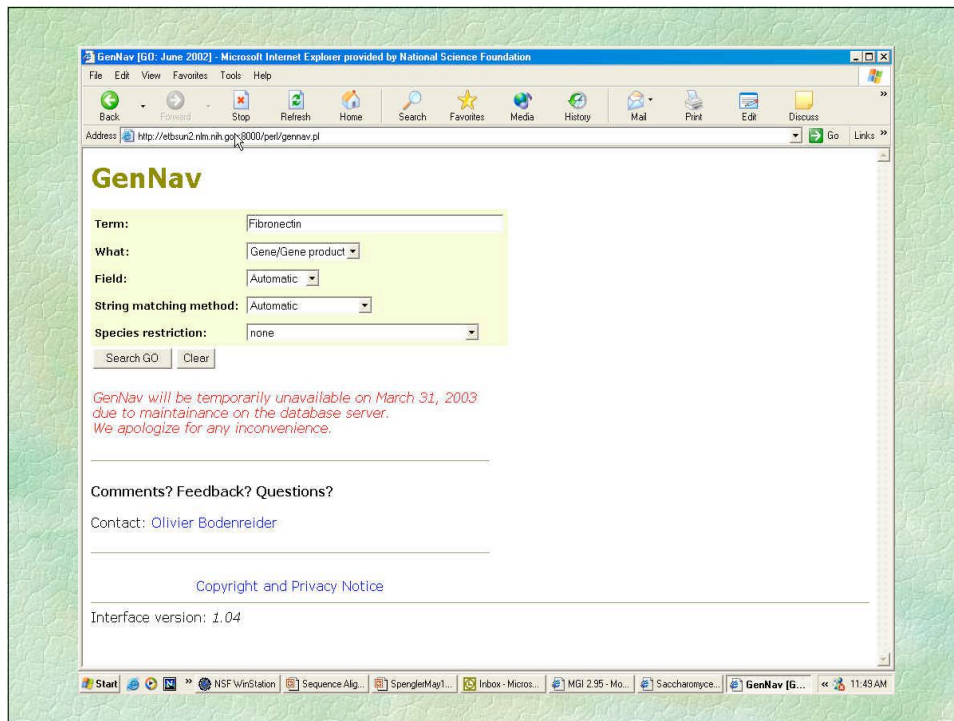
ADDITIONAL INFORMATION for RPE1

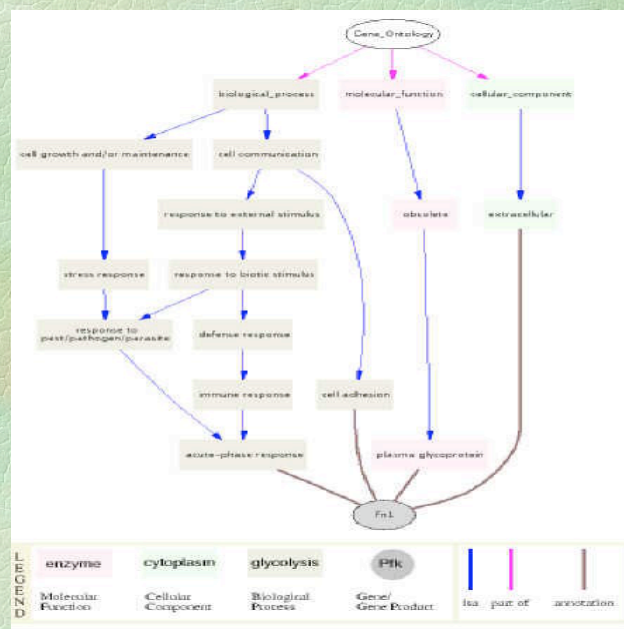
Locus History	Global Gene Hunter	Function Junction	Expression Connection
Protein Info & Composition	Motifs PDB Homologs	Gene/Sequence Resources	

SGD® Pages Database Copyright © 1997-2002 The Board of Trustees of Leland Stanford Junior University. Permission to use the information contained in this database was given to the researcher/institution who contributed or published the information. Users of the database are solely responsible for compliance with any copyright restrictions, including those applying to the author abstracts. Documents from this server are provided "AS-IS" without any warranty, expressed or implied.

[Return to SGD](#) [Send a Message to the SGD Curators](#)



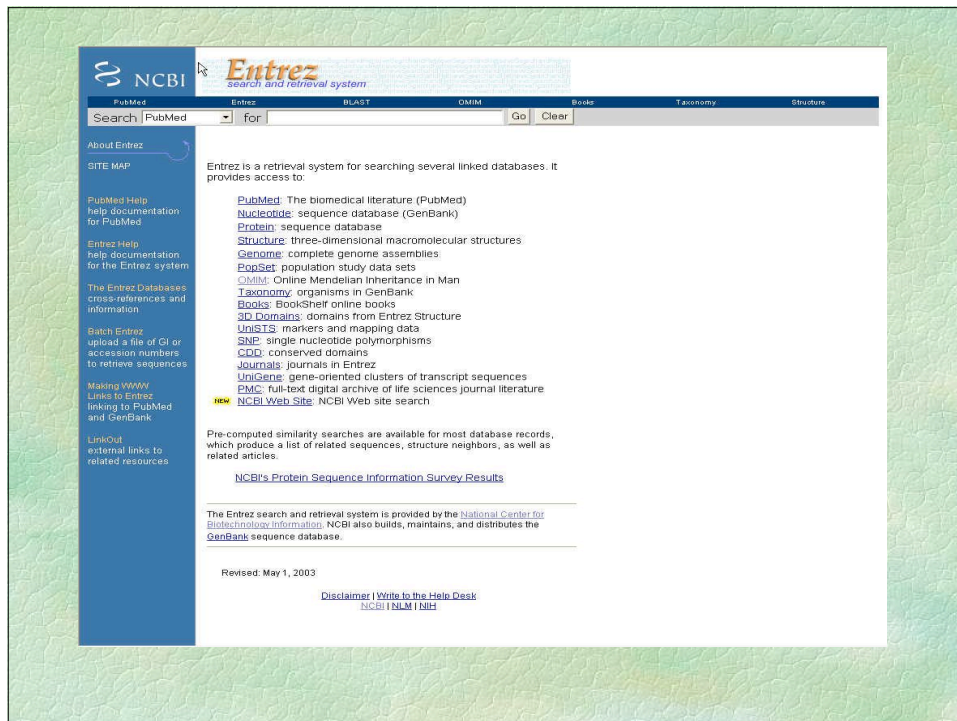




Additional Information

Symbol	Fnl								
Species	0								
Full name	fibronectin 1								
Synonyms	<ul style="list-style-type: none">• Fnl								
Cross-references	<table><tr><th>DB TYPE</th><th>ID</th><th>OBJ_TYPE</th><th>TAXON</th></tr><tr><td>mgc acc</td><td>MG1:95566</td><td>gene</td><td>10090</td></tr></table>	DB TYPE	ID	OBJ_TYPE	TAXON	mgc acc	MG1:95566	gene	10090
DB TYPE	ID	OBJ_TYPE	TAXON						
mgc acc	MG1:95566	gene	10090						
Molecular functions	<ul style="list-style-type: none">• plasma glycoprotein [IEA]								
Biological processes	<ul style="list-style-type: none">• cell adhesion [IEA]• acute-phase response [IEA]								
Cellular components	<ul style="list-style-type: none">• extracellular [IEA]								

New query



ENTREZ: Linked Databases

<http://www.ncbi.nlm.nih.gov/Entrez/>

- Concept of Neighbor-usually BLAST relationship
- Precomputed=Fast
- Related sequence, structure neighbors, related articles
- CUBBY

EST Pluses

- Rapid
 - Inexpensive
 - Applicable to gene discovery, regulation
 - Gene sequence diversity
- Worked to identify over half the human genes

EST negatives

- No defined protein product
- High error rate
- Not curated for annotation
- Low data quality for assembly

EST Data Quality Issues

- Regions of high and low quality reads in the same EST
- Kinds of error
 - Clone orientation
 - Chimeras
 - Missing reads
 - Compression and base calling

Stringent Clustering

- One-pass assembly
- Fewer, shorter consensus sequences
- Lower coverage of expressed gene data
- Lower inclusion rate of expressed gene forms
- Algorithms such as TIGR_ASSEMBLER

Loose Clustering

- One-pass assembly
- Large, “sloppy” clusters
- Greater coverage
- Possible inclusion of paralogous genes
- Lower fidelity
- Includes alternate expressed forms
- Unigene

Steps in EST Clustering

1. Mask for repeats and vector, leaving a minimum number of residues of ‘clean’ data (100bp for Unigene)
2. Initial cluster based on sequence identity
3. Generation of consensus
4. Joining of clusters

STEPS IN CLUSERING-2

☛ Clone joining

- Utilizes the physically shared clone id between 3' and 5' EST fragments from the same clone
 - Can be a source of error because it relies on the accuracy of the annotation and uniqueness of the clone ID, especially when data from disparate sources

What are we trying to do with EST Clustering?

☛ “Match” sequences-either they do or don't

☛ Look at near-identical matches

Tools for Matching

- Smith-Waterman, BLAST, FASTA built for searching-measure *quantitatively* the similarity between any two distances.
- But ESTs either match or not, so need only a near or perfect match.

Adjusting BLAST for ESTs

- Stringent match set for EST
 - E expectation value set to 0
 - G cost to open a gap increased
 - E cost to extend a gap increased
 - Q mismatch penalty increased
 - R match reward increased
 - W word size set for longer words

EST DATABASES:Quality issues

• SEQUENCE QUALITY

- calculated error less than 1% (Phred-20) is the rule
- frameshifts and stops common
- Rules are usually observed by exception
- There are lots of exceptions in the public data
- Many 3' UTRs

EST Databases: Quality #2

• CLONE QUALITY

- Over-representation
- Tissue specificity
- Developmental stage specificity
- Unprocessed mRNA clones
- Chimeras
- Contamination

EST Cluster Databases

- STACK-at SANBI <http://sanbi.ac.za>
- TIGR-animals, plants, other
<http://www.tigr.org/tdb/tgi.shtml>
- Unigene-NCBI
 - Many species, plant, animal,
 - mRNAs
 - predicted mRNAs

UniGene - Microsoft Internet Explorer provided by National Science Foundation

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss Messenger

Address <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=unigene> Go Links »

NCBI UniGene

Search UniGene for Go Clear

Limits Preview/Index History Clipboard Details

NCBI

UniGene is an experimental system for automatically partitioning GenBank sequences into a non-redundant set of gene-oriented clusters. Each UniGene cluster contains sequences that represent a unique gene, as well as related information such as the tissue types in which the gene has been expressed and map location.

Species	Entries
Chordata	
Mammalia	
<i>Bos taurus</i>	12,786 entries
<i>Homo sapiens</i>	108,944 entries
<i>Mus musculus</i>	90,444 entries
<i>Rattus norvegicus</i>	63,253 entries
<i>Sus scrofa</i>	14,390 entries
Aves	
<i>Gallus gallus</i>	5,042 entries
Amphibia	
<i>Silurana tropicalis</i>	7,674 entries
<i>Xenopus laevis</i>	19,045 entries
Actinopterygii	
<i>Danio rerio</i>	15,707 entries
<i>Oryzias latipes</i>	6,115 entries
Ascidacea	
<i>Ciona intestinalis</i>	13,328 entries
Arthropoda	
Insecta	
<i>Anopheles gambiae</i>	3,243 entries
<i>Drosophila melanogaster</i>	14,701 entries
Nematoda	
Chromadorea	

Start NSF WinStation UniGene - Microsoft... Microsoft PowerPoint - [...] Inbox - Microsoft Outlook RE: Great Software - Me... 2:20 PM

NCBI UniGene

PubMed Entrez BLAST OMIM Taxonomy Structure

Search Human Go

NCBI

UniGene

Home Page

Frequently Asked Questions

Query Tips

DDD-Library Digital Differential Display

Download UniGene

UniGene Homo sapiens

Home Page

Release Statistics

Library Report

Library Browser

DDD-Library Digital Differential Display

UniGene Human Sequences Collection

Search using text

At the top of most UniGene pages you will see a text field for use in formulating queries. Query terms may be ordinary words like "kinase" or specific identifiers, such as GenBank accession numbers. In addition, a number of @functions are provided for specialized purposes.

For more information, see our [Query Tips](#)

Search by Chromosome

[1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [8](#) [9](#) [10](#) [11](#) [12](#) [13](#) [14](#) [15](#) [16](#) [17](#) [18](#) [19](#) [20](#) [21](#) [22](#)
[X](#) [Y](#) [Xort](#)

Information on map position has been imported from several sources, including OMIM, the RH Consortium Human Transcript Map, and the Whitehead Institute human physical map. Click on a chromosome number (above) to see a list of UniGene entries for that chromosome. Alternatively, you may use @chr(num) as a search term.

Search by Library

Use the [Library Browser](#) to see a list of cDNA libraries that have been used in EST projects. Note that each library has an ID number, which may also be used within a text query as @lib(num).

[NLM](#) | [NIH](#) | [UniGene](#) | [Privacy Statement](#) | [Disclaimer](#) | [NCBI Help](#)

UNIGENE

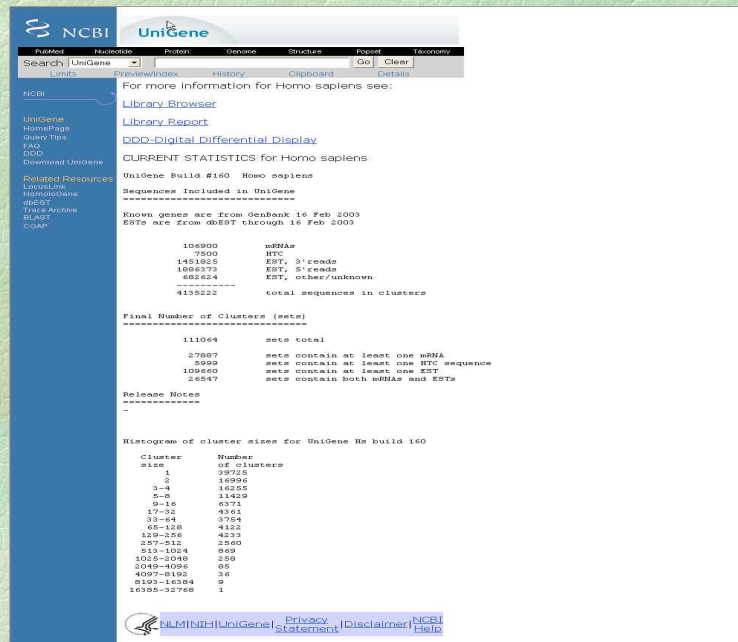
A LIST OF LISTS

- The cluster and known EST, mRNA pieces
- Additional annotation-gene name, etc.
- Distributed as a subset of dbest

NOT included in the BLAST searchable DB at NCBI

Caveats on Clusters

- Not stable
- Can go to complete cDNAs as available



File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address http://www.ncbi.nlm.nih.gov/UniGene/ddd.cgi?ACT=sk%3A20RIG+HsPool1=myeloma+1%3A7038Pool2=New+pool%3ADESC+Hepatic&LID=9631 Go Links

NCBI UniGene

PubMed Nucleotide Protein Genome Structure Popset Taxonomy

Search UniGene Go Clear

Limits Preview/Index History Clipboard Details

Digital Differential Display (DDD)

Begin DDD for:

DDD is a computational method for comparing sequence-based gene representation profiles among individual cDNA libraries or pools of libraries.

[MORE](#)

For the current analysis, the following table describes the pools that have been defined so far.

Pool	Name	Lib ID(s)	Clustered ESTs
Edit...	A	myeloma 1	7038 1802
Edit...	B	Hepatic	9631 13617
New...			

[Start Over](#)

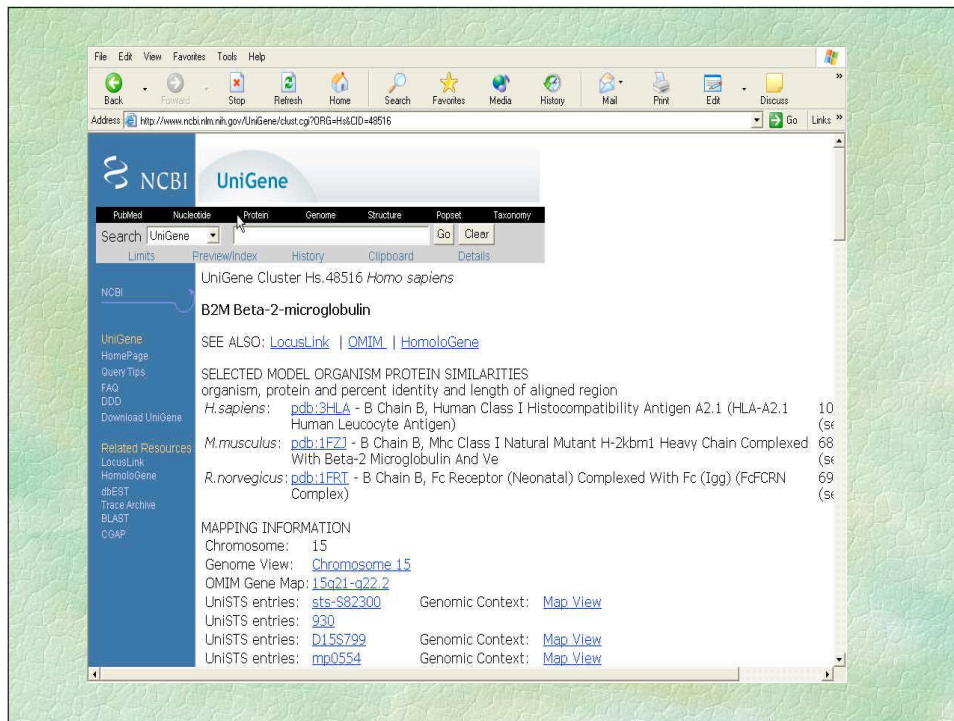
Statistically Significant Differences

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address http://www.ncbi.nlm.nih.gov/UniGene/ddd.cgi?ACT=sk%3A20RIG+HsPool1=myeloma+1%3A7038Pool2=New+pool%3ADESC+Hepatic&LID=9631 Go Links

	0.00000	0.01395	190	Hs.254105 enolase 1, (alpha) (ENO1)
A < B		B > A		
25	0.01387	0.00007	1	Hs.48516 beta-2-microglobulin (B2M)
A > B		B < A		
3	0.00166	0.01278	174	Hs.14376 actin, gamma 1 (ACTG1)
A < B		B > A		
2	0.00111	0.01190	162	Hs.74335 heat shock 90kDa protein 1, beta (HSPCB)
A < B		B > A		
17	0.00943	0.00029	4	Hs.275865 ribosomal protein S18 (RPS18)
A > B		B < A		
15	0.00632	0.00000		Hs.413826 immunoglobulin heavy constant gamma 3 (G3m marker) (IGHG3)
A > B		B < A		
17	0.00943	0.00140	19	Hs.412900 ribosomal protein L10 (RPL10)
A > B		B < A		
	0.00000	0.00720	98	Hs.272499 dehydrogenase/reductase (SDR family) member 2 (DHRS2)
A < B		B > A		
12			12	Hs.401448 tumor protein,



LOCUSLINK

(<http://www.ncbi.nlm.nih.gov/LocusLink>)

- A useful, searchable compendium of loci across *Caenorhabditis elegans*, cow, fruit fly, human, human immunodeficiency virus type 1, mouse, rat, and zebrafish.
- Linked for PubMed, OMIM, RefSeq, Homologene data, Unigene, and Variation Data

NCBI LocusLink

Search: LocusLink Display Brief Organism: All

Query: Go Clear

View: Hs BAT2 One of 1 Loci Save All Loci

ABCDEFGHIJKLMNOPQRSTUVWXYZ

Click to Display mRNA-Genomic Alignments (spanning 16941 bps)

PUB DMM ACEVIEW UNIGENE MAP VAR HOMOL GDB

Homo sapiens Official Gene Symbol and Name (HGNC)

BAT2: HLA-B associated transcript 2

LocusID: 7916

Overview

RefSeq Summary: A cluster of genes, BAT1-BAT5, has been localized in the vicinity of the genes for THF alpha and THF beta. These genes are all within the human major histocompatibility complex class III region. This gene has microsatellite repeats which are associated with the age-at-onset of insulin-dependent diabetes mellitus (IDDM) and possibly thought to be involved with the inflammatory process of pancreatic beta-cell destruction during the development of IDDM. This gene is also a candidate gene for the development of rheumatoid arthritis. There are two alternatively spliced transcripts encoding different isoforms described for this gene.

Protein Summary: Major histocompatibility complex-associated protein; contains proline-rich domains

Locus Type: gene with protein product, function known or inferred

Product: HLA-B associated transcript-2 isoform a
HLA-B associated transcript-2 isoform b

Alternate Symbols: O2, D6S51, D6S51E

Alias: HLA-B associated transcript-2
large proline-rich protein BAT2

Function Submit GeneRIF (All Pubs) ?

Gene Ontology:

Term	Evidence	Source	Pub
MHC-interacting protein	P	Proteome	pm

Relationships

Mouse Homology Maps:

NCBI vs. MGD	UCSC vs. MGD	RefSeq	Protein
17 19.04 cM	Bat2	Hs Mm	
17 19.04 cM	3110039B02Fuk	Hs Mm	

Map Information

Chromosome: 6

Cytogenetic: 6p21.3

Markers: Chr. 6: SHOC
Chr. 7: 34807
B167814

NCBI Reference Sequences (RefSeq)

Category: REVIEWED

1. mRNA: NM_004638

Protein: NP_004628 HLA-B associated transcript-2 isoform b

Transcript Variant: This variant (2) utilizes alternative splice sites in the coding region. It lacks 36 bases, as compared to variant 1 but maintains the same reading frame. Thus isoform b is 12 aa shorter than isoform a.

GenBank Source: AF127256

2. mRNA: NM_006066

Protein: NP_542417 HLA-B associated transcript-2 isoform a

Transcript Variant: This variant (1) encodes the full length isoform.

GenBank Source: AF129756

Category: NCBI Genome Annotation

Genomic Context: NT_007292

Associated transcripts/isoforms for this locus:

gb sw mv ov mm

LocusLink Report - Microsoft Internet Explorer provided by National Science Foundation

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss Links

Address: http://www.ncbi.nlm.nih.gov/LocusLink/LocRpt.cgi?l=567

NCBI LocusLink

PubMed Entrez BLAST OMIM Map Viewer Taxonomy Structure

Search: LocusLink Display Brief Organism: All

Query: Go Clear

View: Hs B2M One of 1 Loci Save All Loci

ABCDEFGHIJKLMNOPQRSTUVWXYZ

Click to Display mRNA-Genomic Alignments (spanning 6627 bps)

PUB DMM ACEVIEW UNIGENE MAP VAR HOMOL GDB

Homo sapiens Official Gene Symbol and Name (HGNC)

B2M: beta-2-microglobulin

LocusID: 567

Overview

Locus Type: gene with protein product, function known or inferred

Product: beta-2-microglobulin

Function Submit GeneRIF (All Pubs) ?

Phenotype: Hemodialysis-related amyloidosis

GeneRIF: Gene References into Function:

11676539 • structure in amyloid fibril formation

11914379 • basis in developing CD8+ T cell

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address <http://www.ncbi.nlm.nih.gov/LocusLink/LocRpt.cgi?l=567> Go Links

LocusLink Home

B2M Index
 Top of Page
 Nomenclature
 Overview
 Function
 Relationships
 Map
 RefSeq
 Related Seqs
 Links

LocusLink
 Collaborators
 Download
 FAQ
 Help
 Statistics

RefSeq
 About
 Download
 FAQ
 Statistics

GeneRIF: Gene References into Function:

- [11676539](#) • structure in amyloid fibril formation
- [11914579](#) • basis in developing CD8+ T cell antagonists
- [12119416](#) • Crystal structure reveals clues to its amyloidogenic properties
- [11967567](#) • Mapping the core of the beta(2)-microglobulin amyloid fibril by H/D exchange
- [11801591](#) • cleaved form partially attains a conformation that has amyloidogenic features
- [11847272](#) • solution structure determined by (1)H NMR spectroscopy and restrained modeling calculations
- [11849381](#) • Signal transduction of beta2m-induced expression of VCAM-1 and COX-2 in synovial fibroblasts.
- [12454016](#) • A dominant negative mutant of this protein blocks the extracellular folding of MHC I heavy chain
- [11943769](#) • Conformation of beta 2-microglobulin

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address <http://www.ncbi.nlm.nih.gov/LocusLink/LocRpt.cgi?l=567> Go Links

LocusLink Home

B2M Index
 Top of Page
 Nomenclature
 Overview
 Function
 Relationships
 Map
 RefSeq
 Related Seqs
 Links

LocusLink
 Collaborators
 Download
 FAQ
 Help
 Statistics

RefSeq
 About
 Download
 FAQ
 Statistics

Gene Ontology™:

Term	Evidence	Source	Pub
• amyloid protein	TAS	GOA	pm
• antigen presentation, endogenous antigen	IEA	GOA	
• antigen processing, endogenous antigen via MHC class I	IEA	GOA	
• MHC class I receptor activity	IEA	GOA	

Relationships ?

Mouse Homology Maps:

NCBI vs. MGD 2 69.00 cM [B2m](#) Hs Mm

Map Information ?

Chromosome: 15 mv

Cytogenetic: 15q21-q22.2 HUGO

Markers:

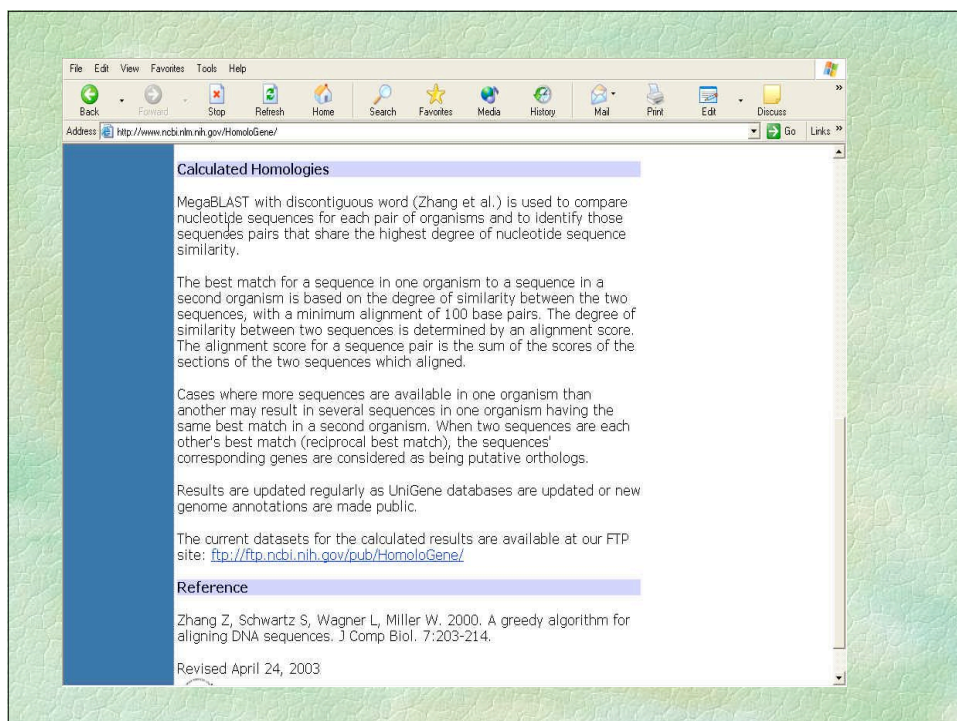
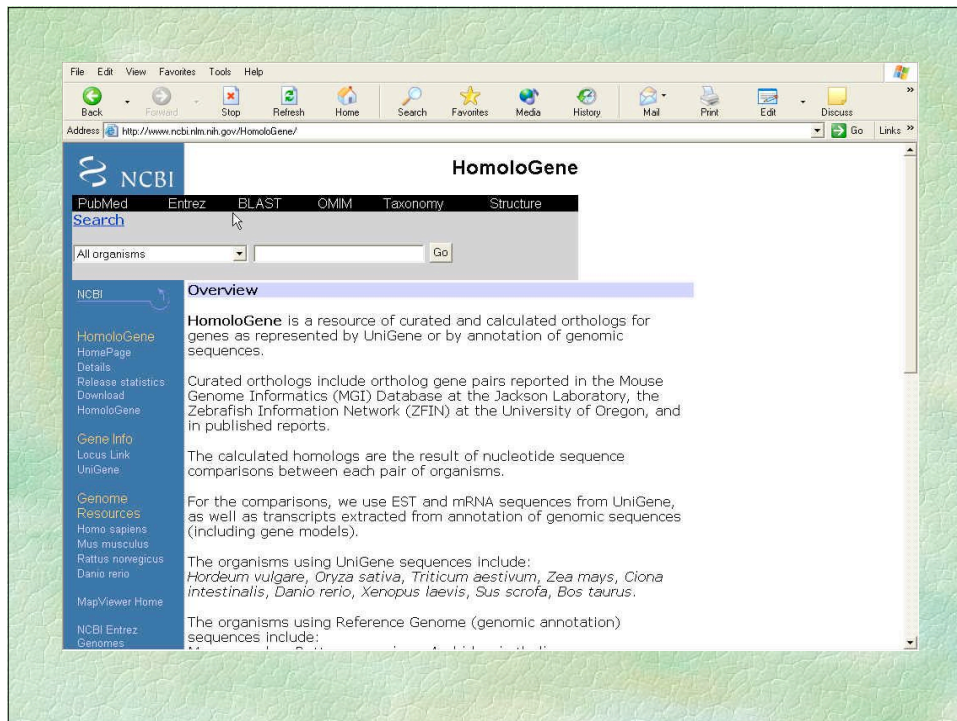
Chr.	Marker	mv
Chr. 15	D15S799	mv
Chr. 15	STS	mv
Chr. 15	S82300	mv
Chr. 15	RH76	mv
Chr. 15	G62079	mv

NCBI Reference Sequences (RefSeq) ?

Category: **PROVISIONAL**

mRNA: [NM_004048](#)

Protein: [NP_004039](#) beta-2-microglobulin BL



File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address <http://www.ncbi.nlm.nih.gov/HomoloGene/homolog?HID=139587> Go Links

NCBI HomoloGene

PubMed Entrez BLAST OMIM Taxonomy Structure

Search

All organisms Go

HOMOLOGENE ENTRY

H.sapiens -B2M beta-2-microglobulin
[UniGene](#) | [LocusLink](#) | [MIM](#) | [MapView](#) | [NM_004048.1](#)

POSSIBLE HOMOLOGOUS GENES

R.norvegicus -B2m Beta-2-microglobulin
[UniGene](#) | [LocusLink](#) | [RGD](#) | [MapView](#) | [NM_012512.1](#)

S.scrofa -Ss.12348 Sus scrofa beta 2-microglobulin mRNA, complete cds.
[UniGene](#) | [U13854.1](#)

B.taurus -B2M beta-2, microglobulin
[UniGene](#) | [LocusLink](#) | [X69084.1](#)

CALCULATED ORTHOLOGS

Listed below are the nucleotide sequence comparisons used in determining homology. The pairs below represent reciprocal best hits; each alignment is the best one for both organisms. The percent ID below represents identity over an aligned region. When present, red arrows (*) point out a group of sequence matches which are part of a triplet, being consistent between more than two organisms.

Organism- Gene	Organism- Gene	Percent ID
H.sapiens -B2M	S.scrofa - Ss.12348	81.0
H.sapiens -B2M	B.taurus - B2M	80.0
H.sapiens -B2M	R.norvegicus - B2m	73.5

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorites Media History Mail Print Edit Discuss

Address <http://www.ncbi.nlm.nih.gov/HomoloGene/homolog?HID=139587> Go Links

CALCULATED ORTHOLOGS

Listed below are the nucleotide sequence comparisons used in determining homology. The pairs below represent reciprocal best hits; each alignment is the best one for both organisms. The percent ID below represents identity over an aligned region. When present, red arrows (*) point out a group of sequence matches which are part of a triplet, being consistent between more than two organisms.

Organism- Gene	Organism- Gene	Percent ID
*H.sapiens -B2M	S.scrofa - Ss.12348	81.0
*H.sapiens -B2M	B.taurus - B2M	80.0
*H.sapiens -B2M	R.norvegicus - B2m	73.5

ADDITIONAL CALCULATED ORTHOLOGS

*S.scrofa -Ss.12348	B.taurus - B2M	85.1
*R.norvegicus -B2m	S.scrofa - Ss.12348	80.3
*R.norvegicus -B2m	B.taurus - B2M	76.7

CURATED ORTHOLOGS

Published orthologs as reported in curated databases

H.sapiens -B2M	R.norvegicus - B2m	PUB
H.sapiens -B2M	M.musculus - B2m	MGI

FURTHER READING

Trinh CH, et al., Crystal structure of monomeric human beta-2-microglobulin reveals clues to its amyloidogenic properties. *Proc Natl Acad Sci U S A* 99, 9771-9776 (2002).

Avet-Loiseau H, et al., Oncogenesis of multiple myeloma: 14q32 and 13q chromosomal abnormalities are not randomly distributed, but correlate with natural history, immunological features, and clinical presentation. *Blood* 99, 2185-2191

Resources for Genomic Comparison

- GLASS-<http://plover.lcs.mit.edu>
- PipMaker: <http://bio.cse.psu.edu>
- Rosetta: [http:// plover.lcs.mit.edu/genes](http://plover.lcs.mit.edu/genes)
- SGP: <http://soft.ice.mpg.de/sgp-1>
- VISTA: <http://www-gsd.lbl.gov/VISTA>
- WABA:
<http://www.cse.ucsc.edu/~kent/xenoAli/index.html>

EFFICIENT TEXT SEARCHING

- Use Wild Cards: #,\$,?,*
- Use Boolean Operators
 - Not
 - And
 - Or
 - Nor

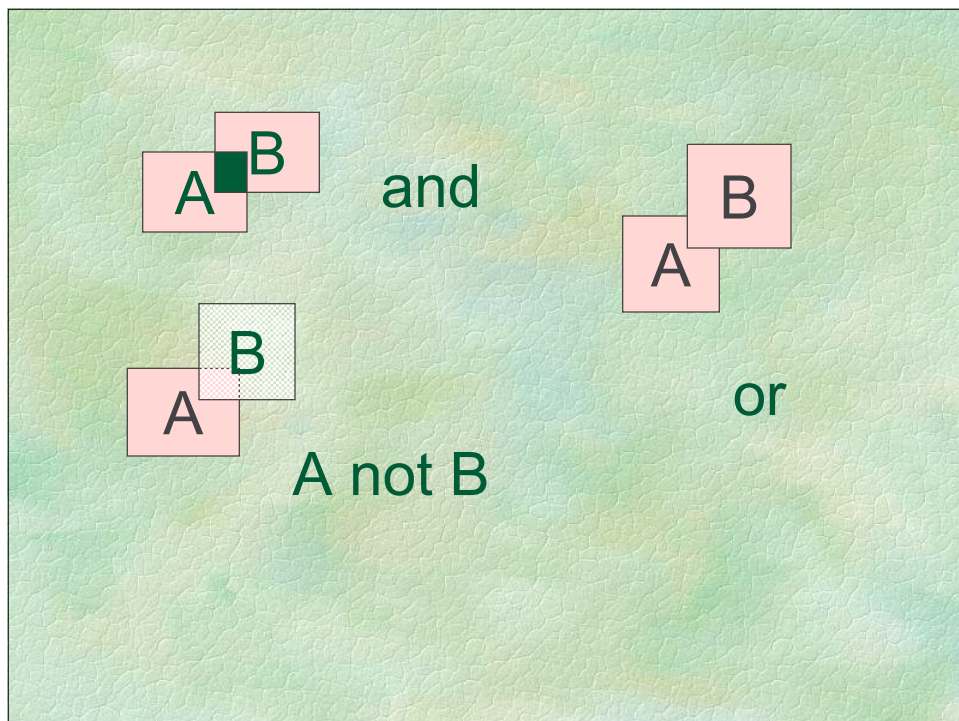
Boolean Operators

🐼 **AND** A and B BOTH

🐼 **OR** A or B EITHER

🐼 **NOT** B not A Have B, do not have A

🐼 **NOR** A nor B A but not B OR B but not A



WILD CARDS

- ☞ Match one character-NCBI uses #
- ☞ Match zero or one character NCBI uses \$, others ?
- ☞ Match zero or more characters-usually *

MEDICAL SUBJECT HEADINGS

- ☞ CONTROLLED Vocabulary
- ☞ Indexing of articles, books, etc.
- ☞ Current version has over 300,000 terms
- ☞ Can download list and make your own assortment

MeSH Advantages

- ☛ Assigned to the the entire document, not just title and abstract
 - ☛ Major topic (*)
 - ☛ Subheadings if available
 - ☛ MeSH topics are exploded to include all the terms included in the meaning.
- Try it; you may like it.

Other Resources

- ☛ NCBI Education Page
<http://www.ncbi.nlm.nih.gov/Education/index.html>
- ☛ BCM Gene Finder
http://searchlauncher.bcm.tmc.edu/docs/sl_links.html
- ☛ EBI-SwissProt, TrEMBL, PIR, SRS, Tools
<http://www.ebi.ac.uk>
- ☛ ExPASy-SwissProt, TrEMBL
<http://www.expasy.ch/>
- ☛ DISC-DNA Information and Stock Center
<http://www.dna.affrc.go.jp>

Final Thoughts

- ☛ Trust your intuition
- ☛ Look at all the possibilities
- ☛ Use all the resources you can

